# Liu-type logistic regression coefficient estimation with multicollinearity using the bootstrapping method

Narumol Sudjai[1] and Monthira Duangsaphon[2*]

[1]*Department of Orthopaedics Surgery, Faculty of Medicine Siriraj Hospital,*
*Mahidol University, Bangkok 10700, Thailand*
[2]*Department of Mathematics and Statistics, Faculty of Science and Technology,*
*Thammasat University, Pathumthani 12121, Thailand*
[*]*Corresponding author: monthira@mathstat.sci.tu.ac.th*

## ABSTRACT

This study proposed new estimators for shrinkage and ridge parameters to overcome the multicollinearity problem in Liu-type logistic regression using the bootstrapping method. Moreover, we compared the performance of four methods for logistic regression coefficient estimation with multicollinearity present: the maximum likelihood estimator, ridge logistic regression, Liu logistic regression, and Liu-type logistic regression, all performed with the bootstrapping method. A simulation study was conducted to compare the performance of the four different estimation methods using the estimated mean square error. The results from both the simulation study and a real data application showed that the Liu-type logistic regression with the bootstrapping method performed best, among the four methods, with a high correlation coefficient. Moreover, the proposed estimators for the shrinkage parameter and ridge parameter showed good performance. In addition, the use of Liu-type logistic regression together using the bootstrapping method was the most robust for correcting the multicollinearity problem.

*Keywords*: bootstrapping method; ridge estimator; Liu estimator; Liu-type estimator; multicollinearity

## 1. INTRODUCTION

Logistic regression analysis is increasingly used in medical and public health research. It's mainly applied to evaluate the relationships between independent and dependent variables, as well as to predict the future value of the dependent variable. For example, in the case where researchers want to study the risk factors for the lymphatic metastasis of malignant bone and soft-tissue tumors, the independent variables of interest include gender, age, location of the primary tumor, and the tumor size, while the dependent variable would simply be coded as 1 (lymph node metastasis) or 0 (non-metastasis). The most commonly used statistical method

in such research is binary logistic regression analysis (Hosmer and Lemeshow, 2000; Petrie and Sabin, 2009; Kleinbaum and Klein, 2010). A critical problem that commonly arises in statistical analysis in medical and public health research is multicollinearity. Multicollinearity is a condition in regression analysis in which some of the independent variables are highly correlated, which can lead to inflating the variance of at least one of the estimated logistic regression coefficients (Petrie and Sabin, 2009; Kleinbaum and Klein, 2010). To solve this problem, we typically exclude some independent variables from the logistic regression model. The choice of which independent variables to select in the model is typically

based on statistical significance, but for some studies, this choice might be difficult if all the independent variables have clinical importance (Petrie and Sabin, 2009). Consequently, several biased estimators have been proposed as alternatives to the maximum likelihood estimator to solve the effect of multicollinearity. For example, Schaefer et al. (1984) proposed a ridge estimator. Later, Urgan and Tez (2008) proposed a Liu estimator, based on the work of Liu (1993). More recently, Huang (2012) proposed a Liu-type estimator involving a combination of two different estimators. Farghali and Abo-El-Hadid (2017) compared the performance of the Liu estimator with that of the maximum likelihood and Stien and ridge estimators, and reported that the Liu estimator was mostly preferred for solving the multicollinearity problem. Also, since the shrinkage parameter (*d*) is an important value for estimating the coefficients in Liu logistic regression, Månsson et al. (2012) proposed five estimators for *d*. Their results showed that their proposed quantiles' shrinkage estimator performed best.

Another problem of using logistic regression analysis in medical and public health research is the typical small sample sizes involved. The bootstrap method is the technique most widely applied to remedy the effect of the small sample sizes (Chernick and La Budde, 2011; Efron, 1979). Sudjai and Duangsaphon (2019) proposed a shrinkage estimator and compared the performance of four methods for performing logistic regression coefficient estimation with the multicollinearity problem present: Liu logistic regression, Liu logistic regression with bootstrapping, almost unbiased Liu logistic regression, and almost unbiased Liu logistic regression with bootstrapping. The results showed that Liu logistic regression with bootstrapping performed best. The shrinkage estimators proposed by Månsson et al. (2012) and Sudjai and Duangsaphon (2019) also performed well. However, the authors did not introduce the exact method for evaluating a single value of *d* and some results showed that the shrinkage estimators are

often equal to zero. Therefore, *d* must be appropriately chosen for each data item.

Furthermore, to the best of our knowledge, no study has yet been performed to compare the performance of Liu-type estimators with the Liu estimator in logistic regression under varying situations and when applied to small sample size data.

Consequently, in this study, we proposed new estimators for the shrinkage and ridge parameters to use with the Liu-type estimator in a logistic regression model using the bootstrapping method under the case with multicollinearity problem present together with small sample sizes. Moreover, the performance of four methods for logistic regression coefficient estimation using the bootstrapping method: the maximum likelihood, ridge, Liu, and Liu-type estimators were compared with the mean square error (MSE) value obtained from Monte Carlo simulations. Furthermore, the applicability of the proposed estimators was demonstrated on a practical data set.

## 2. MATERIALS AND METHODS

Binary logistic regression analysis is used to estimate the logistic regression coefficient. The model is as follows:

$$Y_i = \pi_i + \varepsilon_i \qquad (1)$$

where the dependent variable ($Y_i$) is a dichotomous variable that has a Bernoulli distribution with the parameter value $\pi_i$, where $\pi_i = \frac{e^{x_i\beta}}{1+e^{x_i\beta}}$ ; $i = 1,2,3,\dots,n$. The random error ($\varepsilon_i$) has a distribution with zero mean and a variance equal to $\pi_i(1-\pi_i)$. The independent variable can be both categorical and continuous data (Hosmer and Lemeshow, 2000). We let $X$ be an $n \times (p+1)$ data matrix with $p$ independent variables, where $x_i$ represents the independent variables for the $i^{th}$ row of $X$, $\beta$ is a $(p+1)\times 1$ coefficient vector, and $n$ is the sample size.

The general method used for coefficient estimation in the logistic regression model is the maximum likelihood method. Here, letting $Y_i$ be coded as 1 or 0, the conditional probability that $Y_i$ is equal to one is given by $x_i$, which can be denoted as $P(Y_i = 1|x_i)$. On the other hand, $P(Y_i = 0|x_i)$ is the conditional probability that $Y_i$ is equal to zero, given as $x_i$. For a set of observations $(y_i, x_i)$, if $y_i$ is equal to 1, then the contribution to the likelihood function is $\pi_i$ as well, as if $y_i$ is equal to zero, then the contribution to the likelihood function is $1 - \pi_i$. Thus, the contribution to the likelihood function for the set of observations $(y_i, x_i)$ can be written as:

$$P(Y_i = y_i) = \begin{cases} \pi_i^{y_i}(1 - \pi_i)^{1-y_i} & , y_i = 0,1 \\ 0 & , \text{otherwise.} \end{cases}$$

When the observations are assumed to be independent, the likelihood function can be obtained as follows:

$$l(\underset{\sim}{\beta}) = \prod_{i=1}^{n} \pi_i^{y_i}(1 - \pi_i)^{1-y_i} \qquad (2)$$

and the maximum likelihood estimator can be defined as follows:

$$\hat{\underset{\sim}{\beta}}_{ML} = (X'\widehat{W}X)^{-1}(X'\widehat{W}\hat{\underset{\sim}{z}}) \qquad (3)$$

where $\hat{\underset{\sim}{\beta}}_{ML}$ is a $(p+1) \times 1$ vector of the maximum likelihood estimator; $X$ is an $n \times (p+1)$ data matrix; $\widehat{W}$ is a diagonal matrix of the order $(n \times n)$ with the $i^{th}$ diagonal element equal to $\hat{\pi}_i(1 - \hat{\pi}_i)$ ; $i = 1,2,3,\ldots,n$; and $\hat{\underset{\sim}{z}}$ is an $n \times 1$ vector with the $i^{th}$ element and is defined by:

$$\hat{z}_i = log(\hat{\pi}_i) + \frac{y_i - \hat{\pi}_i}{\hat{\pi}_i(1 - \hat{\pi}_i)}; \ i = 1,2,3,\ldots,n \qquad (4)$$

where $y_i$ is the $i^{th}$ value of the dependent variable and $\hat{\pi}_i$ is the estimator of $\pi_i$.

The mean square error (MSE) of the maximum likelihood estimator can be used to assess the performance of logistic regression coefficient estimators. The MSE of $\hat{\underset{\sim}{\beta}}_{ML}$ can be written as follows:

$$\text{MSE}(\hat{\underset{\sim}{\beta}}_{ML}) = tr[(X'\widehat{W}X)^{-1}] = \sum_{j=1}^{p+1} \frac{1}{\lambda_j} \qquad (5)$$

where $\lambda_j$ is the $j^{th}$ eigenvalue of matrix $X'\widehat{W}X$; $tr[(X'\widehat{W}X)^{-1}]$ is the trace of matrix $(X'\widehat{W}X)^{-1}$; $j = 1,2,3,\ldots, p+1$; and $\hat{\underset{\sim}{\beta}}_{ML}$ is an unbiased estimator. Therefore, $\text{MSE}(\hat{\underset{\sim}{\beta}}_{ML})$ shows the variance of the estimator.

**2.1 Ridge estimator**

In 1984, Schaefer et al. (1984) proposed a ridge estimator, which is a biased estimator, for solving multicollinearity in the logistic regression model. The ridge estimator can be determined as follows:

$$\hat{\underset{\sim}{\beta}}_{RE} = (X'\widehat{W}X + kI)^{-1}(X'\widehat{W}X)\hat{\underset{\sim}{\beta}}_{ML} \qquad (6)$$

where $\hat{\underset{\sim}{\beta}}_{RE}$ is a $(p+1) \times 1$ vector of the ridge estimator and $k$ is the ridge parameter ($k > 0$).

The MSE of the ridge estimator is as follows:

$$\text{MSE}(\hat{\underset{\sim}{\beta}}_{RE}) = \sum_{j=1}^{p+1} \frac{\lambda_j}{(\lambda_j + k)^2} + k^2 \sum_{j=1}^{p+1} \frac{\alpha_j^2}{(\lambda_j + k)^2}$$

$$= f_1(k) + f_2(k) \qquad (7)$$

where $\lambda_j$ is the $j^{th}$ eigenvalue of matrix $X'\widehat{W}X$ and $j = 1,2,3,\ldots,p+1$. Further, $\alpha_j^2$ is the $j^{th}$ element of $\gamma\beta_{ML}$ and $\gamma$ is an eigenvector of matrix $X'\widehat{W}X$. In equation (7), $f_1(k)$ is the variance function and $f_2(k)$ is the square bias.

Several methods can be used for estimating the ridge parameter. However, the choice of the ridge parameter follows no definite rule. In previous studies, various estimators of $k$ have been used in simulation studies (Muniz and Kibria, 2009; Kibria et al., 2012), including the following:

$$\hat{k}_1 = max\left( \frac{1}{\sqrt{\frac{\hat{\sigma}^2}{\hat{\alpha}_j^2}}} \right) \qquad (8)$$

$$\hat{k}_2 = \text{median}\left(\frac{1}{q_j}\right) \tag{9}$$

$$\hat{k}_3 = \prod_{j=1}^{p}\left(\frac{1}{q_j}\right)^{\frac{1}{p}} \tag{10}$$

where $\hat{\sigma}^2 = \frac{\sum_{i=1}^{n}(y_i - \hat{\pi}_i)^2}{n-k-1}$, $q_j = \frac{\lambda_{\max}}{\{(n-p)\hat{\sigma}^2\}+\{\lambda_{\max}\hat{\alpha}_j^2\}}$,

and $\lambda_{max}$ is the maximum eigenvalue of matrix $X'\widehat{W}X$.

## 2.2 Liu estimator

In the 1990s, Liu (1993) proposed the Liu estimator, which is also a biased estimator, for solving the multicollinearity problem, while more recently, Urgan and Tez (2008) proposed an updated Liu estimator. The Liu estimator in logistic regression is defined as follows:

$$\hat{\underset{\sim}{\beta}}_{LE} = (X'\widehat{W}X + I)^{-1}(X'\widehat{W}X + dI)\hat{\underset{\sim}{\beta}}_{ML} \tag{11}$$

where $d$ is the shrinkage parameter $(0 < d < 1)$.

In the case where $d$ equals one, $\hat{\beta}_{LE} = \hat{\beta}_{ML}$. If $d$ is less than one, then $\|\hat{\underset{\sim}{\beta}}_{LE}\| \leq \|\hat{\underset{\sim}{\beta}}_{ML}\|$, where $\|\hat{\underset{\sim}{\beta}}_{LE}\| = \left(\sum_{j=1}^{p+1}\hat{\beta}_{LE,j}^2\right)^{\frac{1}{2}}$ and $\|\hat{\underset{\sim}{\beta}}_{ML}\| = \left(\sum_{j=1}^{p+1}\hat{\beta}_{ML,j}^2\right)^{\frac{1}{2}}$.

The MSE of the Liu estimator is then as follows:

$$\text{MSE}(\hat{\underset{\sim}{\beta}}_{LE}) = \sum_{j=1}^{p+1}\frac{(\lambda_j+d)^2}{\lambda_j(\lambda_j+1)^2} + (d-1)^2\sum_{j=1}^{p+1}\frac{\alpha_j^2}{(\lambda_j+1)^2}$$

$$= f_1(d) + f_2(d) \tag{12}$$

where $f_1(d)$ is the variance function and $f_2(d)$ is the square bias.

Since $d$ is an important value that is commonly applied to remedy multicollinearity, Månsson et al. (2012) proposed five estimators for the shrinkage parameter $(\hat{d})$. These estimators were based on the work of Hoerl and Kennard (1970), Kibria (2003), and Khalaf and Shukur (2005) and could be used to obtain the individual parameter $d_j$ from: $d_j = (\alpha_j^2 - 1)/((1/\lambda_j) + \alpha_j^2)$,

$j = 1,2,3,\ldots,p+1$. This individual parameter can then be used to estimate a single value $\hat{d}$. The results showed that $\hat{d}_1$, which can be considered using quantiles, is the best shrinkage estimator for use in cases involving the presence of multicollinearity. The shrinkage parameter $(\hat{d}_1)$ is given below.

$$\hat{d}_1 = max\left(0, min\left(\frac{\hat{\alpha}_j^2 - 1}{\frac{1}{\hat{\lambda}_j} + \hat{\alpha}_j^2}\right)\right), j = 1,2,3,\ldots,p+1 \tag{13}$$

In 2019, Sudjai and Duangsaphon (2019) proposed two estimators for the shrinkage parameter. The first estimator was based on the work of Hoerl and Kennard (1970), and is as follows:

$$\hat{d}_2 = max\left(0, \frac{\hat{\alpha}_{min}^2 - 1}{\frac{1}{\hat{\lambda}_{min}} + \hat{\alpha}_{min}^2}\right) \tag{14}$$

where $\hat{\alpha}_{min}^2$ is the minimum element of $\gamma\beta_{ML}$, $\gamma$ is the eigenvector of matrix $X'\widehat{W}X$, and $\lambda_{min}$ is the minimum eigenvalue of matrix $X'\widehat{W}X$. As shown in equation (14), $\hat{\alpha}_{min}^2$ and $\hat{\lambda}_{min}$ can be applied to compute $\hat{d}_2$ rather than computing with $\hat{\alpha}_{max}^2$ and $\hat{\lambda}_{max}$ as proposed by Hoerl and Kennard (1970).

The second estimator was based on the work of Månsson et al. (2012), whereby the estimator is computed from:

$$\hat{d}_3 = max\left(0, 1 - \frac{\sum_{j=1}^{p+1}\left\{\frac{1}{\lambda_j(\lambda_j+1)}\right\}}{\sum_{j=1}^{p+1}\left\{\frac{1+\lambda_j\alpha_j^2}{\lambda_j(\lambda_j+1)^2}\right\}}\right), j = 1,2,3,\ldots,p+1 \tag{15}$$

## 2.3 Liu-type estimator

In 2012, Huang (2012) proposed a new Liu-type estimator, which involved a combination of two different estimators and is defined by:

$$\hat{\underset{\sim}{\beta}}_{LTE} = (X'\widehat{W}X + kI)^{-1}(X'\widehat{W}X + kdI)\hat{\underset{\sim}{\beta}}_{ML} \tag{16}$$

where $k > 0$ and $0 < d < 1$. It was shown that if $k = 1$, then $\hat{\beta}_{LTE} = \hat{\beta}_{LE}$; while if $d = 1$ or $k = 0$, then $\hat{\beta}_{LTE} = \hat{\beta}_{ML}$.

The MSE of this Liu-type estimator can be computed from:

$$\text{MSE}(\hat{\beta}_{LTE}) = \sum_{j=1}^{p+1} \frac{(\lambda_j + kd)^2}{\lambda_j(\lambda_j + k)^2} + \sum_{j=1}^{p+1} \frac{k^2(d-1)^2 \alpha_j^2}{(\lambda_j + k)^2}$$

$$= f_1(k, d) + f_2(k, d) \qquad (17)$$

where $f_1(k, d)$ is the variance function and $f_2(k, d)$ is the square bias.

In 2016, Asar (2016) proposed an optimal shrinkage parameter $(\hat{d}_{opt})$, which can be defined by:

$$\hat{d}_{opt} = \frac{\sum_{j=1}^{p+1}\left\{\frac{k\alpha_j^2 - 1}{(\lambda_j + k)^2}\right\}}{\sum_{j=1}^{p+1}\left\{\frac{1 + \lambda_j \alpha_j^2}{\lambda_j(\lambda_j + k)^2}\right\}}, j = 1, 2, 3, \ldots, p+1 \qquad (18)$$

where $-\infty < \hat{d}_{opt} < \infty$ and $k$ is the ridge parameter.

## 2.4 New estimators for the shrinkage and ridge parameters

Considering the previous studies, Schaefer et al. (1984) proposed that $\hat{\beta}_{RE}$ will have a smaller MSE than $\hat{\beta}_{ML}$ for some ridge parameters $(k)$ when the data are collinear and the sample size is sufficiently large. Recently, Urgan and Tez (2008) proposed that $\text{MSE}(\hat{\beta}_{LE}) < \text{MSE}(\hat{\beta}_{ML})$ when the shrinkage parameter is fixed, $0 < d < 1$, and the sample size is sufficiently large. Furthermore, Huang (2012) proposed that $\hat{\beta}_{LTE}$ will have a smaller MSE than $\hat{\beta}_{ML}$ when $k > 0$ and $0 < d < 1$. Therefore, we can see that both $\hat{\beta}_{LE}$ and $\hat{\beta}_{LTE}$ will perform better than $\hat{\beta}_{ML}$ when $d \in (0,1)$. Sudjai and Duangsaphon (2019) performed a simulation study that showed that $\hat{\beta}_{LE}$ performs well when the estimated value of $d$ is close to zero. Moreover, the proposed estimators of $d$ in Månsson et al. (2012) as well as in Sudjai and Duangsaphon (2019) perform well, but have a limitation in that the shrinkage estimators often end up equaling zero based on a 1,000 times simulation.

In this study, we propose a new estimator of $d$. This estimator was developed based on the work of Sudjai and Duangsaphon (2019). Here, to show the optimal parameter $d$, we take the first derivative $\text{MSE}(\hat{\beta}_{LE})$ with respect to $d$ and then equate the derivative to zero. Solving the equation for $d$, we can obtain the estimated value of $d$, which is always less than one. For some points, $d$ may be a negative value because $\lambda_j > 0$. Hence, we replace $\alpha_j^2$ with $(\hat{\alpha}_{ML}^2)_j - \frac{1}{\lambda_j}$, and then letting $d$ have a value between zero and one, we obtain:

$$\hat{d}_4 = max\left( 0, 1 - \frac{\sum_{j=1}^{p+1}\left\{\frac{1}{\lambda_j(\lambda_j+1)}\right\}}{\sum_{j=1}^{p+1}\left\{\frac{1 + \left(\lambda_j\left(\alpha_j^2 - (1/\lambda_j)\right)\right)}{\lambda_j(\lambda_j+1)^2}\right\}} \right), 0 < \hat{d}_4 < 1 \qquad (19)$$

where $\lambda_j$ is the $j^{th}$ eigenvalue of matrix $X'\widehat{W}X$ and $j = 1, 2, 3, \ldots, p+1$; $\alpha_j^2$ is the $j^{th}$ element of $\gamma\beta_{ML}$; and $\gamma$ is the eigenvector of matrix $X'\widehat{W}X$.

After determining $\hat{d}_{opt}$ and $\hat{d}_4$ for $d$, the parameter $k$ must be selected. In this study, we propose a new estimator of $k$ based on the work of Asar (2016). It is easy to find the optimal parameter $k$ by differentiating $\text{MSE}(\hat{\beta}_{LTE})$ with respect to $k$ and then by equating the derivative to zero. Solving the equation for $k$, we obtain the individual parameter as $k_j = \frac{\lambda_j}{\lambda_j\alpha_j^2 - d(\lambda_j\alpha_j^2 + 1)}, \hat{k}_j > 0$, and $j = 1, 2, 3, \ldots, p+1$.

This individual parameter $(k_j)$ was used to estimate a single value $(\hat{k}_{LTE})$, based on the works of Hoerl and Kennard (1970) and Månsson et al. (2012). Then, $\hat{k}_{LTE}$ is as follows:

$$\hat{k}_{LTE} = max\left(0, \frac{\hat{\lambda}_{min}}{\hat{\lambda}_{min}\hat{\alpha}^2_{min} - \hat{d}(\hat{\lambda}_{min}\hat{\alpha}^2_{min}+1)}\right) \qquad (20)$$

where $\hat{\alpha}^2_{min}$ is defined as the minimum element of $\gamma\beta_{ML}$ and $\gamma$ is the eigenvector of matrix $X'\hat{W}X$. Also, $\lambda_{min}$ is defined as the minimum eigenvalue of matrix $X'\hat{W}X$. We substitute the unknown parameters with $\hat{\alpha}^2_{min}$ and $\hat{\lambda}_{min}$ in line with ideas taken from Hoerl and Kennard (1970) and Sudjai and Duangsaphon (2019).

## 2.5 Bootstrapping method

Multicollinearity and a small sample size are key problems in logistic regression coefficient estimation, and can lead to a poor performance of the estimators (Kleinbaum and Klein, 2010; Stoltzfus, 2011). To remedy the effects of these problems, we constructed biased estimators, namely, the maximum likelihood, ridge, Liu, and Liu-type estimators, under the bootstrapping technique. There are two approaches for bootstrapping: the first approach involves resampling the random error term ($\varepsilon_i$), while the second approach resamples from the observations. In this study, we chose the second approach. Hence, the proposed process for bootstrapping is as follows.

Step 1. Creating a bootstrap sample of size $n$ ($z_1^*, z_2^*, \ldots, z_n^*$) from the original data with the replacement giving $\frac{1}{n}$ probability for each $z_i^*$. Thus, we can obtain the following: $z_i^* = (y_i^*, x_i^*)$, $i = 1,2,3,\ldots,n$.

Step 2. Estimating $\beta$ for the logistic regression model using the four different estimation methods: maximum likelihood method; Liu logistic regression together with the shrinkage parameters $\hat{d}_1$, $\hat{d}_2$, and $\hat{d}_3$; ridge logistic regression together with the ridge parameters $\hat{k}_1$, $\hat{k}_2$, and $\hat{k}_3$; and Liu-type logistic regression together with $\hat{k}_1\hat{d}_{opt}$ and $\hat{k}_{LTE}\hat{d}_4$.

Step 3. Repeating steps 1 and 2 for $B$ times, where $B$ is the number of repetitions. We can thus obtain bootstrap estimates of the parameter $\beta$ for each estimator.

Step 4. Using the resulting bootstrap estimates in step 3 (e.g., $\hat{\beta}_{ML}^{*(1)}, \hat{\beta}_{ML}^{*(2)}, \ldots, \hat{\beta}_{ML}^{*(B)}$) to compute the

average estimate for each parameter. We can thus obtain the estimated value of the parameters for use with the bootstrapping method as follows:

$$\bar{\hat{\beta}}_{ML}^* = \frac{1}{B}\sum_{b=1}^{B}\hat{\beta}_{ML}^{*(b)} \qquad (21),$$

$$\bar{\hat{\beta}}_{RE(\hat{k}_1)}^* = \frac{1}{B}\sum_{b=1}^{B}\hat{\beta}_{RE(\hat{k}_1)}^{*(b)} \qquad (22),$$

$$\bar{\hat{\beta}}_{RE(\hat{k}_2)}^* = \frac{1}{B}\sum_{b=1}^{B}\hat{\beta}_{RE(\hat{k}_2)}^{*(b)} \qquad (23),$$

$$\bar{\hat{\beta}}_{RE(\hat{k}_3)}^* = \frac{1}{B}\sum_{b=1}^{B}\hat{\beta}_{RE(\hat{k}_3)}^{*(b)} \qquad (24),$$

$$\bar{\hat{\beta}}_{LE(\hat{d}_1)}^* = \frac{1}{B}\sum_{b=1}^{B}\hat{\beta}_{LE(\hat{d}_1)}^{*(b)} \qquad (25),$$

$$\bar{\hat{\beta}}_{LE(\hat{d}_2)}^* = \frac{1}{B}\sum_{b=1}^{B}\hat{\beta}_{LE(\hat{d}_2)}^{*(b)} \qquad (26),$$

$$\bar{\hat{\beta}}_{LE(\hat{d}_3)}^* = \frac{1}{B}\sum_{b=1}^{B}\hat{\beta}_{LE(\hat{d}_3)}^{*(b)} \qquad (27),$$

$$\bar{\hat{\beta}}_{LTE(\hat{k}_1,\hat{d}_{opt})}^* = \frac{1}{B}\sum_{b=1}^{B}\hat{\beta}_{LTE(\hat{k}_1,\hat{d}_{opt})}^{*(b)} \qquad (28),$$

$$\bar{\hat{\beta}}_{LTE(\hat{k}_{LTE},\hat{d}_4)}^* = \frac{1}{B}\sum_{b=1}^{B}\hat{\beta}_{LTE(\hat{k}_{LTE},\hat{d}_4)}^{*(b)} \qquad (29).$$

where $\hat{\beta}_{ML}^{*(b)}$ is the bootstrap estimates of the maximum likelihood estimator,

$\hat{\beta}_{RE(\hat{k})}^{*(b)}$ is the bootstrap estimates of the ridge estimator,

$\hat{\beta}_{LE(\hat{d})}^{*(b)}$ is the bootstrap estimates of the Liu estimator,

$\hat{\beta}_{LTE(\hat{k},\hat{d})}^{*(b)}$ is the bootstrap estimates of the Liu-type estimator and $b = 1,2,3,\ldots,B$.

## 2.6 Monte Carlo simulation

The key factors that can affect the performance of the estimation methods are the number of independent variables, the sample size, and the degree

of correlation among the independent variables. These factors were thus varied in the simulation study.

In this study, Monte Carlo simulations were performed using 2 and 4 independent variables ($p$) (Månsson et al., 2012). The sample size ($n$) was equal to 50, 100, 200, and 400. For $n < 30p$, the sample size ($n$) is regarded as a small sample case (Kerlinger and Pedhazur, 1973). Therefore, the situation was either $n = 50$ when $p = 2$ or $n = 50$ or 100 when $p = 4$, referring to a small sample case. The degree of correlation ($\rho$) was varied at 0.1, 0.3, 0.5, 0.75, 0.85, and 0.95, which represent positive correlations. The independent variables were generated from:

$$x_{ij} = (1 - \rho^2)^{\frac{1}{2}} z_{ij} + \rho z_{ip} \, ,$$
$$i = 1,2,3,...,n, \text{ and } j = 1,2,3,...,p \qquad (30)$$

where $z_{ij}$ are pseudo-random numbers generated from the standard normal distribution. The dependent variable was generated from a Bernoulli distribution with the parameter $\pi_i$, where $\pi_i = \frac{e^{x_i\beta}}{1+e^{x_i\beta}}$. The values of parameter $\beta$ were chosen so that , $\beta'\beta = 1$ where $\beta_0 = 0$ and $\beta_1, \beta_2, \ldots, \beta_p = \frac{1}{\sqrt{p}}$ based on Newhouse and Oman (1971), who explained that if the MSE is a function of the parameters $\beta$, $\sigma^2$, and $k$ and the number of independent variables are fixed, then the MSE is minimized when we choose the coefficient vector ($\beta'\beta = 1$). Moreover, the intercept value ($\beta_0$) is another important factor because it equals the average value of the log odds ratio. So, if $\beta_0 = 0$, then there is an equal average probability of obtaining one and zero (Månsson and Shukur, 2011). A simulation study was conducted to compare the performance of the four different estimation methods with the estimated MSE. In this study, the experiment was repeated 1,000 times by generating new original data based on Inan and Erdogan (2013), while the bootstrap replication ($B$) was performed 500 times based on Davidson and Mac Kinnon (2000) and Pattengale et al. (2010). The estimated MSE could be computed from:

$$\text{MSE} = \frac{1}{1,000} \sum_{i=1}^{1,000} (\hat{\beta} - \beta)' (\hat{\beta} - \beta) \qquad (31)$$

where $\hat{\beta}$ is a ($p$+1)×1 vector of the logistic regression coefficient estimator and $i$ is the experiment number, in which $i = 1,2,3,...,1,000$.

The estimation method with the lowest estimated MSE is considered the best option for solving the multicollinearity problem in the logistic regression model.

## 3. RESULTS AND DISCUSSION
### 3.1 Simulation study

The estimated MSE values for all the methods when $p = 2$ and 4, $\rho = 0.1$, 0.3, 0.5, 0.75, 0.85, and 0.95, and for different $n$ are listed in Tables 1 and 2, respectively. The effect of increasing $n$ while holding $\rho$ and $p$ fixed was most commonly a decrease in the estimated MSE. While with an increase in the $\rho$ level, the estimated MSE values of all the methods increased when $n$ and $p$ were fixed. In the cases of $\rho = 0.1$, 0.3, and 0.5 in Table 1, we found that the estimated MSE values of the maximum likelihood estimator were less than for the ridge, Liu, and Liu-type logistic regressions with the bootstrapping method. For $\rho = 0.75$, 0.85, and 0.95 in Table 1, the inflation of the estimated MSE values for the Liu-type logistic regression was less than for the maximum likelihood estimator and ridge logistic regression with the bootstrapping method for some $k$, and for the Liu logistic regression with the bootstrapping method. However, there were some differences between the performance of the Liu-type logistic regression with the bootstrapping method, which depended on the estimators of $d$ and $k$. We can see that $\hat{k}_{LTE}\hat{d}_4$ was the best option for all situations. Comparing the performance of this method with the other methods, Liu-type logistic regression with $\hat{k}_{LTE}\hat{d}_4$ based on the bootstrapping method was mostly preferred for correcting the multicollinearity with small sample sizes ($n < 30p$) problem in the logistic regression model.

**Table 1** Estimated MSE values for different $k$ and $d$, when $p = 2$

| $\rho$ | $n$ | BML | BRE | | | BLE | | | BLTE | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | $\hat{k}_1$ | $\hat{k}_2$ | $\hat{k}_3$ | $\hat{d}_1$ | $\hat{d}_2$ | $\hat{d}_3$ | $\hat{k}_1\hat{d}_{opt}$ | $\hat{k}_{LTE}\hat{d}_4$ |
| 0.10 | 50 | 0.4171 | 0.4584 | 0.4312 | 0.4323 | 0.4327 | 0.4327 | 0.4324 | 0.4498 | 0.7335 |
| | 100 | 0.3606 | 0.3727 | 0.3651 | 0.3657 | 0.3658 | 0.3658 | 0.3658 | 0.3715 | 0.6618 |
| | 200 | 0.3502 | 0.3582 | 0.3548 | 0.3550 | 0.3553 | 0.3553 | 0.3553 | 0.3581 | 0.6044 |
| | 400 | 0.3408 | 0.3441 | 0.3427 | 0.3428 | 0.3429 | 0.3429 | 0.3429 | 0.3441 | 0.5503 |
| 0.30 | 50 | 0.5462 | 0.6125 | 0.5832 | 0.5853 | 0.5859 | 0.5859 | 0.5858 | 0.6092 | 0.8730 |
| | 100 | 0.4570 | 0.4739 | 0.4661 | 0.4664 | 0.4674 | 0.4674 | 0.4674 | 0.4735 | 0.7259 |
| | 200 | 0.3988 | 0.4061 | 0.4025 | 0.4028 | 0.4033 | 0.4033 | 0.4033 | 0.4060 | 0.6867 |
| | 400 | 0.3875 | 0.3902 | 0.3892 | 0.3893 | 0.3894 | 0.3894 | 0.3894 | 0.3902 | 0.6198 |
| 0.50 | 50 | 0.7586 | 0.7775 | 0.7628 | 0.7632 | 0.7689 | 0.7689 | 0.7691 | 0.7779 | 0.9390 |
| | 100 | 0.4944 | 0.5135 | 0.5036 | 0.5042 | 0.5068 | 0.5068 | 0.5068 | 0.5134 | 0.7969 |
| | 200 | 0.4718 | 0.4830 | 0.4787 | 0.4790 | 0.4808 | 0.4808 | 0.4808 | 0.4830 | 0.7873 |
| | 400 | 0.4380 | 0.4408 | 0.4397 | 0.4398 | 0.4400 | 0.4400 | 0.4400 | 0.4408 | 0.5864 |
| 0.75 | 50 | 1.1713 | 1.0398 | 1.0749 | 1.0697 | 1.0798 | 1.0798 | 1.0816 | 1.0484 | 1.0074 |
| | 100 | 0.8552 | 0.6629 | 0.7397 | 0.7356 | 0.7509 | 0.7509 | 0.7560 | 0.7016 | 0.6068 |
| | 200 | 0.6173 | 0.5946 | 0.6033 | 0.6027 | 0.6044 | 0.6044 | 0.6045 | 0.5968 | 0.5741 |
| | 400 | 0.6035 | 0.5886 | 0.5952 | 0.5951 | 0.5945 | 0.5945 | 0.5945 | 0.5892 | 0.5595 |
| 0.85 | 50 | 1.8986 | 1.3368 | 1.4321 | 1.4128 | 1.5239 | 1.5237 | 1.5337 | 1.3847 | 1.1279 |
| | 100 | 1.3660 | 0.8922 | 1.0322 | 1.0254 | 1.0833 | 1.0833 | 1.0893 | 0.9372 | 0.8244 |
| | 200 | 0.8082 | 0.7520 | 0.7750 | 0.7736 | 0.7783 | 0.7783 | 0.7792 | 0.7622 | 0.6795 |
| | 400 | 0.7188 | 0.7029 | 0.7094 | 0.7090 | 0.7094 | 0.7094 | 0.7094 | 0.7037 | 0.6663 |
| 0.95 | 50 | 2.0969 | 1.5144 | 1.5559 | 1.5432 | 1.6948 | 1.6859 | 1.7662 | 1.6691 | 1.4612 |
| | 100 | 1.5876 | 1.2492 | 1.2712 | 1.2694 | 1.3353 | 1.3344 | 1.3520 | 1.2972 | 1.1481 |
| | 200 | 1.1911 | 0.8498 | 0.9081 | 0.9101 | 0.9655 | 0.9654 | 0.9699 | 0.8745 | 0.8442 |
| | 400 | 0.7987 | 0.7669 | 0.7754 | 0.7765 | 0.7795 | 0.7795 | 0.7795 | 0.7666 | 0.7320 |

Note: BML = maximum likelihood estimator with the bootstrapping method, BRE = Ridge logistic regression with the bootstrapping method, BLE = Liu logistic regression with the bootstrapping method, and BLTE = Liu-type logistic regression with the bootstrapping method

From Table 2, in the case of $p = 4$, $\rho = 0.1$ and 0.3, the inflation of the estimated MSE values with the ridge logistic regression with the $\hat{k}_1$ bootstrapping method was less than for the maximum likelihood estimator, and Liu and Liu-type logistic regressions with the bootstrapping method. For $\rho = 0.5, 0.75, 0.85,$ and 0.95, the inflation of the estimated MSE values with the Liu-type logistic regression with the bootstrapping method was less than for the maximum likelihood estimator and ridge logistic regression with the bootstrapping method for some $k$, and for the Liu logistic regression with the bootstrapping method. However, there were differences between the performances of Liu-type logistic regression with the bootstrapping method, which depended on the estimators of $d$ and $k$. We can see that $\hat{k}_{LTE}\hat{d}_4$ was the best option for all situations. Thus, Liu-type logistic regression with $\hat{k}_{LTE}\hat{d}_4$ based on the bootstrapping method under the multicollinearity problem with small sample sizes showed the best performance when compared with the other methods. For $n = 50$, $\rho = 0.95$, and $p = 4$, the estimated MSE values for the ridge logistic regression with $\hat{k}_1$ based on the

bootstrapping method was close to the estimated MSE values for Liu-type logistic regression with $\hat{k}_{LTE}\hat{d}_4$ based on the bootstrapping method. Furthermore, increasing the number of independent variables caused an increase in the estimated MSE values of all methods when $\rho$ and $n$ were fixed.

**Table 2** Estimated MSE values for different $k$ and $d$, when $p = 4$

| $\rho$ | $n$ | BML | BRE | | | BLE | | | BLTE | |
|--------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| | | | $\hat{k}_1$ | $\hat{k}_2$ | $\hat{k}_3$ | $\hat{d}_1$ | $\hat{d}_2$ | $\hat{d}_3$ | $\hat{k}_1\hat{d}_{opt}$ | $\hat{k}_{LTE}\hat{d}_4$ |
| 0.10 | 50 | 0.7856 | 0.6079 | 0.6958 | 0.6899 | 0.6793 | 0.6793 | 0.6809 | 0.6197 | 0.7910 |
| | 100 | 0.6586 | 0.6049 | 0.6302 | 0.6282 | 0.6271 | 0.6271 | 0.6271 | 0.6062 | 0.7697 |
| | 200 | 0.5084 | 0.4845 | 0.4965 | 0.4956 | 0.4952 | 0.4952 | 0.4952 | 0.4848 | 0.6508 |
| | 400 | 0.4591 | 0.4557 | 0.4571 | 0.4570 | 0.4569 | 0.4569 | 0.4569 | 0.4558 | 0.6630 |
| 0.30 | 50 | 1.0807 | 0.9283 | 0.9900 | 0.9847 | 0.9863 | 0.9863 | 0.9870 | 0.9356 | 0.9450 |
| | 100 | 0.8502 | 0.7889 | 0.8201 | 0.8181 | 0.8149 | 0.8149 | 0.8149 | 0.7910 | 0.8180 |
| | 200 | 0.6630 | 0.6467 | 0.6547 | 0.6541 | 0.6528 | 0.6528 | 0.6528 | 0.6468 | 0.7910 |
| | 400 | 0.4787 | 0.4765 | 0.4773 | 0.4773 | 0.4771 | 0.4771 | 0.4771 | 0.4766 | 0.7026 |
| 0.50 | 50 | 1.2371 | 0.9942 | 1.0933 | 1.0839 | 1.0864 | 1.0864 | 1.0896 | 1.0151 | 0.9740 |
| | 100 | 1.0522 | 0.8396 | 0.9496 | 0.9420 | 0.9340 | 0.9340 | 0.9341 | 0.8438 | 0.8333 |
| | 200 | 0.8717 | 0.7899 | 0.8388 | 0.8363 | 0.8295 | 0.8295 | 0.8295 | 0.7946 | 0.7542 |
| | 400 | 0.7417 | 0.7255 | 0.7350 | 0.7346 | 0.7323 | 0.7323 | 0.7323 | 0.7256 | 0.6930 |
| 0.75 | 50 | 1.6975 | 1.1999 | 1.3392 | 1.3267 | 1.3791 | 1.3790 | 1.3976 | 1.2387 | 1.0208 |
| | 100 | 1.3885 | 0.9287 | 1.1645 | 1.1553 | 1.1417 | 1.1417 | 1.1436 | 0.9627 | 0.8458 |
| | 200 | 1.0563 | 0.8653 | 0.9718 | 0.9643 | 0.9533 | 0.9533 | 0.9542 | 0.8853 | 0.8159 |
| | 400 | 0.8601 | 0.7977 | 0.8381 | 0.8356 | 0.8272 | 0.8272 | 0.8272 | 0.8014 | 0.7581 |
| 0.85 | 50 | 2.0154 | 1.4542 | 1.6291 | 1.6085 | 1.7199 | 1.7198 | 1.7567 | 1.5346 | 1.1342 |
| | 100 | 1.5301 | 1.1101 | 1.2890 | 1.2753 | 1.2885 | 1.2885 | 1.2992 | 1.1628 | 0.9817 |
| | 200 | 1.0653 | 0.9434 | 1.0016 | 0.9989 | 0.9932 | 0.9932 | 0.9940 | 0.9519 | 0.9410 |
| | 400 | 0.9342 | 0.8971 | 0.9165 | 0.9154 | 0.9120 | 0.9120 | 0.9120 | 0.8974 | 0.8432 |
| 0.95 | 50 | 4.6803 | 1.9512 | 2.2069 | 2.1792 | 2.7193 | 2.7071 | 3.0694 | 2.3924 | 1.9509 |
| | 100 | 2.1355 | 1.2544 | 1.4107 | 1.3946 | 1.5416 | 1.5412 | 1.6231 | 1.3689 | 1.1484 |
| | 200 | 1.9152 | 1.1665 | 1.3980 | 1.3724 | 1.4589 | 1.4589 | 1.4959 | 1.2685 | 1.0406 |
| | 400 | 1.1655 | 0.9516 | 1.0478 | 1.0421 | 1.0438 | 1.0438 | 1.0493 | 0.9785 | 0.8455 |

Note: BML = maximum likelihood estimator with the bootstrapping method, BRE = Ridge logistic regression with the bootstrapping method, BLE = Liu logistic regression with the bootstrapping method, and BLTE = Liu-type logistic regression with the bootstrapping method

## 3.2 Real data application

In this section, a real data application was presented to show the performance when using the new estimators for the shrinkage and ridge parameters to estimate the Liu-type estimator in the logistic regression model using the bootstrapping method under the case with the multicollinearity problem present and small sample sizes. A data set from the UCI machine learning repository (https://archive.ics.uci.edu/ml/index.php) was used fifty donors at random from the data set were selected and then the data were modeled using a binary logistic regression model. The dependent variable was coded as 1 if a blood sample from a donor had good quality blood components and 0 if the sample did not

have good quality blood components. The independent variables were the following:

X₁: frequency (total number of donations),

X₂: time (months since first donation).

The correlation between X₁ and X₂ was equal to 0.96 (high correlation coefficient).

From Table 3, we estimated the standard errors of the different estimators obtained by the bootstrapping method. When looking at the standard errors, we can see that the lowest standard errors were obtained with the Liu-type logistic regression together with $\hat{k}_{LTE}\hat{d}_4$, while the largest were obtained with the maximum likelihood estimator. This means that the Liu-type logistic regression together with $\hat{k}_{LTE}\hat{d}_4$ using the bootstrapping method showed better performance than the other methods for solving the multicollinearity problem with small sample sizes in the logistic regression model, which corresponds to the results obtained from this simulation study.

**Table 3** Estimated parameters ($\hat{\beta}$) and standard errors (*se*) for each method

| Variables | BML | BRE | | | BLE | | | BLTE | |
|---|---|---|---|---|---|---|---|---|---|
| | | $\hat{k}_1$ | $\hat{k}_2$ | $\hat{k}_3$ | $\hat{d}_1$ | $\hat{d}_2$ | $\hat{d}_3$ | $\hat{k}_1\hat{d}_{opt}$ | $\hat{k}_{LTE}\hat{d}_4$ |
| X₁ | | | | | | | | | |
| $\hat{\beta}_1$ | 0.7135 | 0.5000 | 0.6813 | 0.6684 | 0.5982 | 0.5982 | 0.6460 | 0.6385 | 0.5385 |
| $se(\hat{\beta}_1)$ | 0.0091 | 0.0067 | 0.0081 | 0.0075 | 0.0074 | 0.0074 | 0.0088 | 0.0087 | 0.0062 |
| X₂ | | | | | | | | | |
| $\hat{\beta}_2$ | -0.1494 | -0.1116 | -0.1438 | -0.1416 | -0.1292 | -0.1292 | -0.1374 | -0.1357 | -0.1124 |
| $se(\hat{\beta}_2)$ | 0.0019 | 0.0017 | 0.0018 | 0.0017 | 0.0018 | 0.0018 | 0.0019 | 0.0019 | 0.0012 |

Note: BML = maximum likelihood estimator with the bootstrapping method, BRE = Ridge logistic regression with the bootstrapping method, BLE = Liu logistic regression with the bootstrapping method, and BLTE = Liu-type logistic regression with the bootstrapping method

From the simulated MSE values in Tables 1 and 2, it can be seen that the factors influencing the estimated MSE values were: the correlation coefficient level ($\rho$), the sample size ($n$), and the number of independent variables ($p$). An increase in the correlation coefficient level led to an increase in the estimated MSE values for all methods, while holding $p$ and $n$ fixed. The worst case was obtained when the sample size was small ($n = 50$) and $\rho$ was high ($\rho = 0.95$). In the case when $\rho$ and $p$ were fixed, when the sample size increased, then the estimated MSE values of all the methods decreased in all situations. Moreover, an increase in the number of independent variables caused an increase in the estimated MSE values for all methods, while holding $n$ and $\rho$ fixed. Thus, if the number of independent variables is increased, one should also increase the sample sizes to obtain stable estimates. In this study, the method with the best performance was the Liu-type logistic regression together with $\hat{k}_{LTE}\hat{d}_4$ using the bootstrapping method for solving multicollinearity in the logistic regression model. However, if the data are very highly collinear and the sample size is very small, then the performance of ridge logistic regression with $\hat{k}_1$ is close to that of Liu-type logistic regression together with $\hat{k}_{LTE}\hat{d}_4$.

Finally, the estimation methods were applied to a real data set, where the effect of changing the total number of donations and months since the first donation on the quality of a blood sample was explored, and it was shown that the Liu-type logistic regression with $\hat{k}_{LTE}\hat{d}_4$ using the bootstrapping method was the best method.

## 4. CONCLUSION

The MSE was estimated for four logistic regression coefficient estimation methods with different $\rho$, $n$, and $p$, and the results revealed that Liu-type logistic regression with the bootstrapping method performed better than Liu logistic regression with the bootstrapping method in the case of a high correlation coefficient. However, there were some situations where ridge logistic regression with the bootstrapping method for some $k$ showed better performance than Liu-type logistic regression with the bootstrapping method, and this depended on the estimators of $d$ and $k$. Thus, the choice of $d$ and $k$ must be appropriate. In the simulation study, the proposed estimators of the shrinkage parameter $\hat{d}_4$ and ridge parameter $\hat{k}_{LTE}$ showed good performance when the data were highly collinear. Moreover, the Liu-type logistic regression together with $\hat{k}_{LTE}\hat{d}_4$ using the bootstrapping method was the most robust for solving the multicollinearity problem. However, if the data are weakly/moderately collinear and $p = 2$, then the performance of the maximum likelihood estimator with the bootstrapping method will be the best option. In addition, if the data are weakly collinear and $p = 4$, then ridge logistic regression with the bootstrapping method will perform best.

## REFERENCES

Asar, Y. (2016). New shrinkage parameters for the Liu-type logistic estimators. *Communication in Statistics - Simulation and Computation*, 45(3), 1094-1103.

Chernick, M. R., and La Budde, R. A. (2014). *An Introduction to Bootstrap Methods with Applications to R*, New Jersey: John Wiley & Sons.

Davidson, R., and Mac Kinnon, J. G. (2000). Bootstrap tests: How many bootstraps? *Econometric Reviews*, 19(1), 55-68.

Efron, B. (1979). Bootstrap methods: another look at the Jackknife. *The Annals of Statistics*, 7(1), 1-26.

Farghali, R. A., and Abo-El-Hadid, S. M. (2017). Evaluating the performance of Liu logistic regression estimator. *Research Journal of Mathematics and Statistics*, 9(2), 11-19.

Hoerl, A. E., and Kennard, R. W. (1970). Ridge regression: biased estimation for nonorthogonal problems. *Technometrics*, 12(1), 55-67.

Hosmer, D. W., and Lemeshow, S. J. (2000). *Applied Logistic Regression*, 2nd, New York: Wiley & Sons.

Huang, J. (2012). A Simulation research on a biased estimator in logistic regression model. In *Computational Intelligence and Intelligent Systems. ISICA 2012. Communications in Computer and Information Science* (Li, Z., Li, X., Liu, Y., and Cai, Z., eds.), pp. 389-395. Berlin, Heidelberg: Springer.

Inan, D., and Erdogan, B. E. (2013). Liu-type logistic estimator. *Communications in Statistics - Simulation and Computation*, 42(7), 1578-1586.

Kerlinger, F. N., and Pedhazur, E. J. (1973). *Multiple Regression in Behavioral Research*, New York: Holt, Rinehart & Winston.

Khalaf, G., and Shukur, G. (2005). Choosing ridge parameter for regression problems. *Communication in Statistics - Theory and Methods*, 34(5), 1177-1182.

Kibria, B. M. G., Månsson, K., and Shukur, G. (2012). Performance of some logistic Ridge regression estimators. *Computational Economics*, 40(4), 401-414.

Kibria, B. M. G. (2003). Performance of some new ridge regression estimators. *Communications in Statistics - Simulation and Computation*, 32(2), 419-435.

Kleinbaum, D. G., and Klein, M. (2010). *Logistic Regression; A Self-learning Text*, 3rd, New York: Springer.

Liu, K. (1993). A new class of biased estimate in linear regression. *Communications in Statistics - Theory and Methods*, 22, 393-402.

Månsson, K., Kibria, B. M. G., and Shukur, G. (2012). On Liu estimators for the logit regression model. *Economic Modelling*, 29(4), 1483-1488.

Månsson, K., and Shukur, G. (2011). On ridge parameters in logistic regression. *Communications in Statistics - Theory and Methods*, 40, 3366-3381.

Muniz, G., and Kibria, B. M. G. (2009). On some ridge regression estimators: an empirical comparisons. *Communication in Statistics - Simulation and Computation*, 38(3), 621-630.

Newhouse, J. P., and Oman, S. D. (1971). *An Evaluation of Ridge Estimators*, USA: Rand Corporation.

Pattengale, N. D., Alipour, M., Bininda-Emonds, O. R., Moret, B. M., and Stamatakis, A. (2010). How many bootstrap replicates are necessary? *Journal of Computational Biology*, 17(3), 337-354.

Petrie, A., and Sabin, C. (2009). *Medical Statistics at a Glance*, 3rd, Oxford: Wiley Blackwell.

Schaefer, R. L., Roi, L. D., and Wolfe, R. A. (1984). A ridge logistic estimator. *Communications in Statistics - Theory and Methods*, 13(1), 99-113.

Stoltzfus, J. C. (2011). Logistic regression: a brief primer. *Academic Emergency Medicine*, 18(10), 1099-1104.

Sudjai, N., and Duangsaphon, M. (2019). Liu logistic regression coefficient estimation with multicollinearity problem by using the bootstrapping method. *Veridian E-Journal, Science and Technology*, 6(1), 47-61.

Urgan, N. N., and Tez, M. (2008). Liu estimator in logistic regression when the data are collinear. In *20th EURO Mini Conference Continuous Optimization and Knowledge Based Technologies*, Neringa, Lithuania.